

Inspiron 15
3000

29

SATURDAY
2022 JANUARYWEEK 05
029-336TutorialProbability Distributions

2021 DECEMBER

M	T	W	T	F	S
5	6	7	8	9	10
12	13	14	15	16	17
19	20	21	22	23	24
26	27	28	29	30	31

Probability Distributions

A statistical model that shows possible outcomes of a particular event or course of action as well as the statistical likelihood of each event.

10

y-axis
Density

11

12

1. How to go about this?

2. How do we use the collected business data (Sales Volume, loan defaulters, Salary Hikes in an organization etc)?

3. The data values themselves are used directly in the simulation. This is called trace-driven simulation

30

2. "Fit" a theoretical distribution to the data (and check whether that "fit" is good!)
3. The data values could be used to define an empirical distribution function in some way.

FEBRUARY 2022

M	T	W	T	F	S	S
1	2	3	4	5	6	
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28						

MONDAY
JANUARY 2022

31

WEEK 06
031-334What are these empirical distributions?

- Using the data, we build our own distributions.
 - How does one build a distribution?
 - Essential building blocks:
- Define the density/distribution functions.
Estimate the parameters (mean, standard deviation, etc).

Empirical distributions

For ungrouped data:

Let $X_{(i)}$ denote the i^{th} smallest of the X_j 's so that: $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$.

$$F(x) = \begin{cases} 0 & \text{if } x < X_{(1)} \\ \frac{i-1}{n-1} + \frac{x - X_{(i)}}{(X_{(i+1)} - X_{(i)})} & \text{if } X_{(i)} \leq x < X_{(i+1)} \\ 1 & \text{if } X_{(n)} \leq x \end{cases}$$

For Grouped data:

$$G(x) = \begin{cases} 0 & \text{if } x < a_0 \\ G(a_{j-1}) + \frac{x - a_{j-1}}{a_j - a_{j-1}} [G(a_j) - G(a_{j-1})] & \text{if } a_{j-1} \leq x < a_j, j=1, 2, \dots, K \\ 1 & \text{if } a_K \leq x \end{cases}$$

1

TUESDAY

2022 FEBRUARY

WEEK 06
032-333

2022 JANUARY

S	M	T	W	T	F	S
30	31					1
2	3	4	5	6	7	8
9	10	11	12	13	14	15
16	17	18	19	20	21	22
23	24	25	26	27	28	29

The three approaches..

- Approach 1 is used to validate simulation model output for an existing system with the corresponding output for the system itself.
- Two drawbacks of approach 1: Simulation can only reproduce only what happened historically; and there is seldom enough data to make all simulation runs.
- Approaches 2 and 3 avoid these shortcomings so that any value b/w minimum and maximum can be generated. So approaches 2 and 3 are preferred over approach 1.
- If theoretical distributions can be found that fits the observed data (approach 2), then it is preferred over approach 3.

2

: static background

6 value b_i

$$x_{i+1} = \frac{1}{b_i} \left[(x_i - a_i) + \frac{1}{b_i} \right] \quad i = 1, 2, \dots, n-1$$

value b_i

$$\frac{1}{b_i}$$

1

Business Example

2

WEDNESDAY

2022 FEBRUARY

2022 JANUARY

S	M	T	W	T	F	S
30	31					1
2	3	4	5	6	7	8
9	10	11	12	13	14	15
16	17	18	19	20	21	22
23	24	25	26	27	28	29

WEEK 06
033-332

Business Example

8. Clues from Summary statistics

- 9. • For the symmetric distribution, mean and median should match. In the sample data, if these values are sufficiently closer to each other, we can think of a symmetric distribution (e.g. normal)
- 10. • Coefficient of variation (cv): (ratio of std. dev and mean = σ/μ) for continuous distributions. The cv=1 for exponential dist. If the histogram look like a slightly right skewed curve with $cv > 1$, then lognormal could be better approximation of the distribution.
- 11. • Lewis ratio: same as cv for discrete distributions.
- 12. • Skewness(v): measure of symmetry of a distribution
- 13. For normal dist. $v=0$. For $v>0$, the distribution is skewed towards right (exponential dist, $v=2$).
14. And for $v<0$, the distribution is skewed towards left. (left tail is longer than right tail)

Parameter estimation

- Once distribution is guessed, the next step is estimating the parameters of the distribution.
- Each distribution has a set of parameters.
- ✓ Normal distribution has mean and standard deviation
- ✓ Exponential distribution has a " λ ".
- Most common method of parameter estimation: MLE (most likelihood estimation).

MARCH 2022

M	T	W	T	F	S	S
1	2	3	4	5	6	
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31			

Checking the distributionTHURSDAY
FEBRUARY 2022

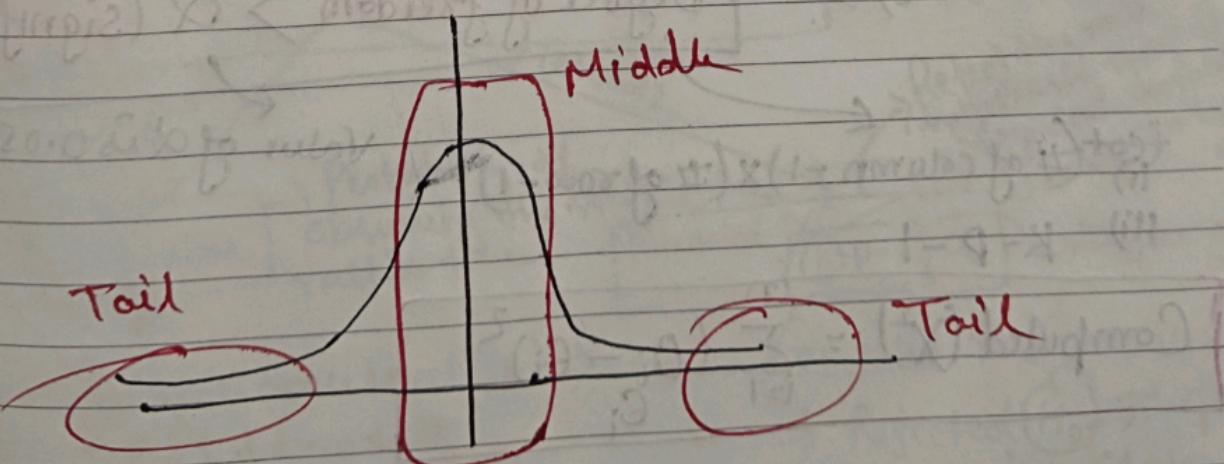
3

WEEK 06
034-331Goodness of fit

- For the input data, we have consumed a probability distribution.
- We also have estimated the parameters for the same.
- It can be checked by several methods:
- 1. Frequency comparison (a bit technical)
- 2. Probability plots (visual tool)
- 3. Goodness-of-fit tests (statistical test of goodness. Very widely used).

Probability plots

- The Q-Q plot will amplify differences between the tails of the model distribution and the sample distribution.
- Whereas, the P-P plot will amplify the differences at the middle portion of the model and sample distribution.



(Date)

4

FRIDAY

2022 FEBRUARY

WEEK 06
035-330

2022 JANUARY

S	M	T	W	T	F	S
30	31					1
2	3	4	5	6	7	8
9	10	11	12	13	14	15
16	17	18	19	20	21	22
23	24	25	26	27	28	29

Goodness-of-fit tests

- A goodness-of-fit test is a statistical hypothesis test that is used to assess formally whether the observations $X_1, X_2, X_3, \dots, X_n$ are an independent sample from a particular distribution with function F .

H_0 : The X_i 's are IID (Independent Identical Distributed) random variables with distribution function F .

- Two famous tests:

- Chi-square test

- Kolmogorov-Smirnov test

Reject the Null Hypothesis

Tabulated (χ^2) =

Degree of freedom > α (significance level)

feet (# of column - 1) x (# of row - 1)

i) K - P - 1

Value of α is 0.05 or 0.01

$$\text{Computed } (\chi^2) = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

NUARY

MARCH 2022

M	T	W	T	F	S	S
1	2	3	4	5	6	
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31			

TutorialSATURDAY
FEBRUARY 2022

5

Overview of Chi-square GOF

WEEK 06
036-329

Fit

Chi-square GOF (Goodness of Fit)

This test is basically going to be used to see if the proposed distribution or hypothesized distribution for the population is a good fit for the population or not.

Frequency Table

Bins	Observed freq.	Expected freq.	χ^2
1	O_1	E_1	$(O_1 - E_1)^2 / E_1$
2	O_2	E_2	$(O_2 - E_2)^2 / E_2$
.	.	.	.
.	.	.	.
.	.	.	.
K	O_K	E_K	$(O_K - E_K)^2 / E_K$

Test statistic: $\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i}$

Method 1
no. of parameters in population that need to compute for given test

Method 2

Probability of observing the sample if null is true

null is true

Reject the null hypothesis that the categorical variables are independent

Tabulated $\chi^2_{(K-p-1), \alpha}$

no. of bins

Degree of freedom

$P\text{-value} < \text{Significance level } (\alpha) \Rightarrow$
Reject H_0 (dependent)

$\chi^2_{\text{Tabulated}} < \chi^2_{\text{Computed}}$

\Rightarrow Rejected H_0 \Rightarrow the population is uniform (independent)

H_a : The prob. is not uniform (not independent)

Inspiron 15
3000 Series

2022 JANUARY

S	M	T	W	T	F	S
30	31					1
2	3	4	5	6	7	8
9	10	11	12	13	14	15
16	17	18	19	20	21	22
23	24	25	26	27	28	29

17

MONDAY

2022 FEBRUARY

WEEK 07
038-327

Population Follows a

K-2-1

(i) Uniform Distribution \Rightarrow two parameters will be calculating \Rightarrow lower value (a), upper value (b) \Rightarrow So for uniform distribution $P = 2(a, b)$

K-1-1

(ii) Poisson distribution \Rightarrow only one parameter (λ) \Rightarrow So $P = 1(\lambda)$

K-2-1

(iii) Normal Distribution \Rightarrow two parameters (μ, σ^2)
 \Rightarrow Two parameters $P = 2(\mu, \sigma^2)$

Tutorial on performing Chi-square GOF for a Uniform Distribution

MARCH

MARCH 2022

M	T	W	T	F	S	S
1	2	3	4	5	6	
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31			

TUESDAY

FEBRUARY 2022

8

WEEK 07

099-16

Step 1:- Descriptive statistics } Narrowing down on possible
Step 2:- Visual tools } guesses (for the parent distribution)

Step 3:- Start creating the frequency table (computing the test statistics)

Step 4:- Conclusion

(H₀) Null Hypothesis:- The given data follows Uniform distribution.

(H_a) Alternate hypothesis:- The given data does not follow Uniform Distribution.

MARCH

APRIL